

## Online Convex Optimization (Convexity, FTL)

Lecturer: Kris Kitani

Scribes: Haowen Shi, Fan Jia

# 1 Review

## 1.1 Overview

In the past lectures we have studied a few online learning algorithms including PWEA and Online Linear Classification. In this lecture we study a generalized framework of these online learning algorithms called *Online Convex Optimization* (OCO). We explain why OCO is a generalization of online learning in section 2.2.

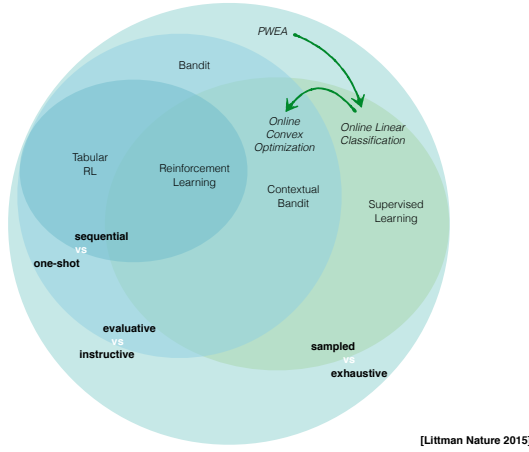


Figure 1: Where we are in the class

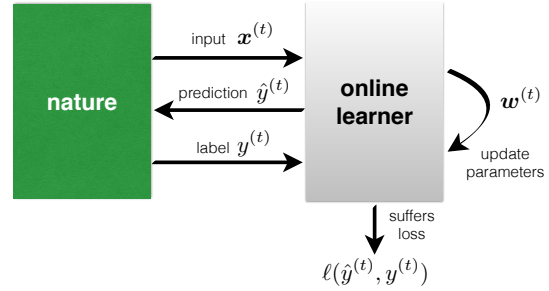


Figure 2: Online learning framework

## 1.2 Online Learning

In online learning, the prediction model is incrementally improved over multiple iterations. The events that happen in each iteration is illustrated in Fig. 2. At each time step  $t$ , the learner receives input  $\mathbf{x}^{(t)}$ , makes a prediction  $\hat{y}^{(t)}$  and then receives the true outcome  $y^{(t)}$  from the environment. The loss function  $\ell(\hat{y}^{(t)}, y^{(t)})$  is evaluated by the learner to produce a loss which drives the update of model weights  $\mathbf{w}^{(t)}$ .

In previous lectures we have learned about two main classes of online learning:

- Prediction with Expert Advice (PWEA)
  - Greedy / Consistent Algorithm
  - Halving Algorithm

- Weighted Majority Algorithm
- Randomized Weighed Majority Algorithm
- Online Linear Classification
  - Perception Algorithm
  - Winnow Algorithm

---

**Algorithm 1** Perceptron algorithm
 

---

```

1:  $\mathbf{w}^{(1)} \leftarrow \mathbf{0}$  ▷ Weight initialization
2: for  $t = 1, \dots, T$  do
3:   RECEIVE  $(\mathbf{x}^{(t)} \in \mathbb{R}^N)$  ▷ Receive expert predictions
4:    $\hat{y}^{(t)} = \text{sign}(\langle \mathbf{w}^{(t)}, \mathbf{x}^{(t)} \rangle)$  ▷ Make learner prediction
5:   RECEIVE  $(y^{(t)} \in \{-1, 1\})$  ▷ Receive actual answer
6:    $\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} + y^{(t)} \cdot \mathbf{x}^{(t)} \cdot \mathbf{1}[y^{(t)} \neq \hat{y}^{(t)}]$  ▷ Weight update
7: end for
  
```

---



---

**Algorithm 2** Winnow algorithm
 

---

```

1:  $\mathbf{w}^{(1)} = \{1, \dots, 1\}$  ▷ Weight initialization
2: for  $t = 1, \dots, T$  do
3:   RECEIVE  $(\mathbf{x}^{(t)} \in \{0, 1\}^N)$ 
4:    $\hat{y}^{(t)} = \mathbf{1}[\langle \mathbf{w}^{(t)}, \mathbf{x}^{(t)} \rangle > N]$  ▷ Make prediction
5:   RECEIVE  $(y^{(t)} \in \{0, 1\})$ 
6:    $w_i^{(t+1)} = w_i^{(t)}(1 + \beta)^{(y^{(t)} - \hat{y}^{(t)}) \cdot x_i^{(t)}}$  ▷ Weight update
7: end for
  
```

---

## 2 Online Convex Optimization

### 2.1 The Optimization Problem

An optimization problem is the process of finding the best solution among all feasible solutions. Given a function of variables  $f(\mathbf{x})$ , we want to find the set of variables  $\mathbf{x}$  such that the function produces the optimal (minimum) value  $\hat{\mathbf{x}}$ .

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} f(\mathbf{x})$$

In real life, most optimization problems are also *constrained* to meet a set of requirements, these constraints can fall in the following categories:

- Inequality Constraint:  $g(\mathbf{x}) \leq 0$
- Equality Constraint:  $h(\mathbf{x}) = 0$
- Domain Constraint:  $\mathbf{x} \in \mathbb{S}$

To solve the optimization problem, there are three high level approaches:

### 1. Analytic solution

In some optimization problems we have a well defined objective function  $f$  so we can find its minima by evaluating its gradient  $\nabla f$  and setting it to zero. Analytic solution is fast and global for convex objective functions but it is not always possible because the objective function might be too complex to evaluate or sometimes hidden to us.

### 2. Brute force search

Brute force search works well in small domain sizes and always obtains globally optimal solutions. The computation soon becomes intractable as the domain size grows so it's not always feasible.

### 3. Numerical methods

Solving optimization problems using numerical methods is very popular in the learning community because the computation cost scales well as problem difficulty grows. However, the trade off is numerical solutions are not always globally optimal because they can converge to local minima depending on the initialization.

## 2.2 Relationship with Online Learning

Online Convex Optimization is a more generalized framework for online learning. In OCO, the learner does not make any observations of the environment, it just sends a set of parameters  $\mathbf{w}^{(t)}$  to the environment and gets a loss function  $l^{(t)}$  back, as shown in Fig. 3 and Algorithm 3. The reason we say OCO is a higher level framework is that we can easily use OCO to solve online learning problems by replacing the "nature" in Fig. 3 with a online learning model. As illustrated in Fig. 4, the online learner in red is calling the online convex optimizer in gray for the prediction step to update its weights.

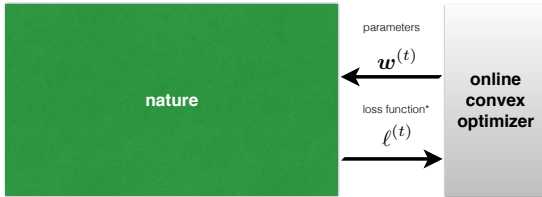


Figure 3: OCO framework

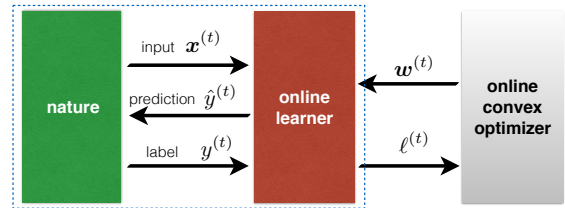


Figure 4: OCO is a generalization

---

### Algorithm 3 Online Convex Optimization

---

```

1: function ONLINECONVOPT(Convex set  $\mathcal{S}$ )
2:   for  $t = 1, 2, \dots, T$  do
3:     PREDICT ( $\mathbf{w}^{(t)} \in \mathcal{S}$ )
4:     RECEIVE ( $f_t : \mathcal{S} \rightarrow \mathbb{R}$ )
5:     COMPUTE LOSS ( $f_t(\mathbf{w}^{(t)})$ )
6:   end for
7: end function

```

---

$\triangleright$  Solution space must be a convex set  
 $\triangleright$  Loss function must be convex

---

## 2.3 Convex Optimization

Generally if a optimization problem is convex, there will be a global optima which can be found faster than non-convex problems.

### 2.3.1 Convex Set

A set  $\mathcal{S}$  is convex if for all  $\mathbf{w}, \mathbf{v} \in \mathcal{S}$ , and for all  $\alpha \in [0, 1]$ , the following holds:

$$\alpha \mathbf{w} + (1 - \alpha) \mathbf{v} \in \mathcal{S}$$

Geometrically, for any pair of points  $\mathbf{w}$  and  $\mathbf{v}$  in  $\mathcal{S}$ , all the points that lie on the line segment between  $\mathbf{w}$  and  $\mathbf{v}$  should also be within the set  $\mathcal{S}$ , as shown in Fig. 5.

### 2.3.2 Convex Function

A function  $f : \mathcal{S} \rightarrow \mathbb{R}$  is convex if for all  $\mathbf{w}, \mathbf{v} \in \mathcal{S}$ , and for all  $\alpha \in [0, 1]$ , the following holds:

$$f(\alpha \mathbf{w} + (1 - \alpha) \mathbf{v}) \leq \alpha f(\mathbf{w}) + (1 - \alpha) f(\mathbf{v})$$

Geometrically, this means for any two points on the function, the secant line between two points always lies above the curve segment connecting the two points. Shown in Fig. 6, a function is convex if the green secant line is always higher than the blue curve segment between the point pair for all pair of points in the function.

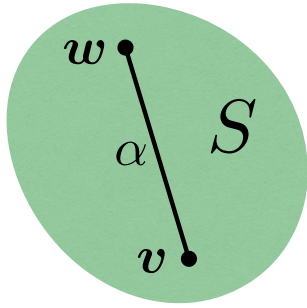


Figure 5: Convex set visualization

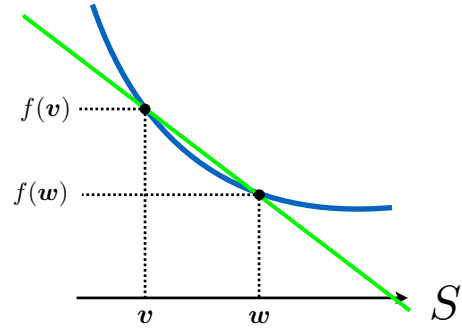


Figure 6: Convex function visualization

Convex optimization requires a convex-Lipschitz-bounded (or convex-smooth-bounded) function and a convex solution space. In the next section, we discuss what is *Lipschitz Continuity*.

### 2.3.3 Lipschitz Continuity

A function  $f(\cdot)$  is called L-Lipschitz smooth over a set  $\mathcal{S}$  with respect to a norm  $\|\cdot\|$  if for all  $\mathbf{u}, \mathbf{w} \in \mathcal{S}$ , we have:

$$|f(\mathbf{u}) - f(\mathbf{w})| \leq L \|\mathbf{u} - \mathbf{w}\| \quad (1)$$

The inequality (1) describes that the rate of change of function  $f$  has to be upper bounded by  $L\|\mathbf{u} - \mathbf{w}\|$  where  $L$  is a Lipschitz constant, which limits how fast the function can change. Fig. 7 shows a good visualization [2] of this upper bound: for a Lipschitz-continuous function there exists a double cone (white) that could be moved along the curve without ever intersecting with the curve itself.

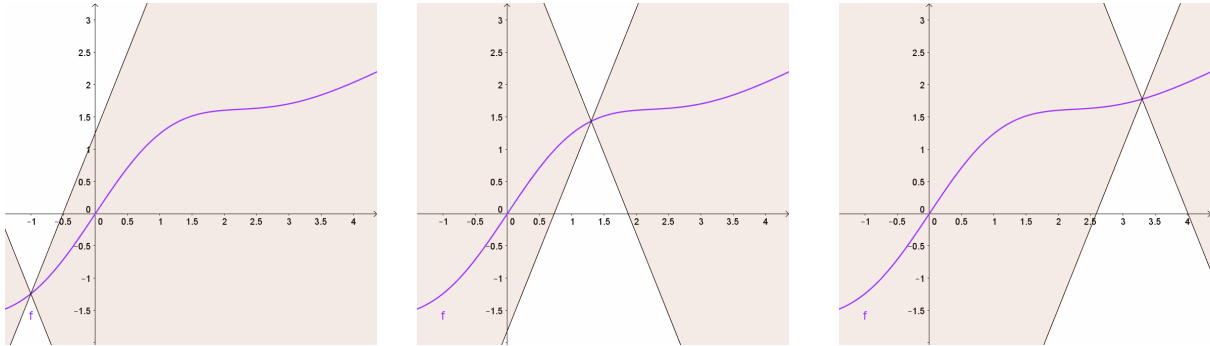


Figure 7: Lipschitz continuity visualization [2]

Note that the Lipschitz continuity has more to do with the analysis of performance guarantee of the algorithm than whether the optimization will work or not. The Lipschitz constant is a mathematical tool for algorithmic analysis instead of a hyperparameter.

### 2.3.4 Convexification

A loss function is easier to solve when it is convex. If a loss function is not convex, we can use some convexification tricks to make it so. In this section we discuss two methods to convert non-convex loss functions to convex ones.

#### a. Convexification by randomization

Recall the weighted majority algorithm we studied before where we make prediction by taking the sign of inner product between observation and weights:

$$\hat{y}^{(t)} = \text{sign} \langle \mathbf{x}^{(t)}, \mathbf{w}^{(t-1)} \rangle$$

We define the loss function as:

$$l^{(t)} = \mathbf{1}[\hat{y}^{(t)} \neq y^{(t)}]$$

This zero-one loss function is non-convex. Fig. 8, shows this with an example pair of points forming a segment that passes beneath the function.

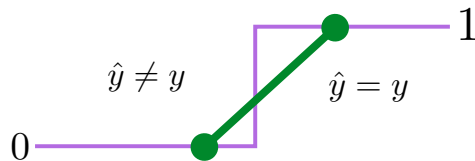


Figure 8: Zero-one loss function is non-convex

Recall that in randomized weighed majority algorithm, we changed the prediction rule to randomly sample an expert's advice from a multinomial distribution.

$$i \sim \text{MULTINOMIAL}(\mathbf{w}^{(t-1)} / \phi^{(t-1)})$$

It turns out that based on such random sampling, the loss function is changed to be an expectation over the zero-one loss.

$$l^{(t)} = E_{\mathbf{p}} \left[ \mathbf{1}[y^{(t)} \neq \hat{y}^{(t)}] \right] = \sum_{p_n} p_n \cdot \mathbf{1}[y^{(t)} \neq \hat{y}_n^{(t)}]$$

Now because of the expectation, the loss function becomes linear and could take any values between zero and one. Therefore, by doing randomization, we converted the non-convex zero-one loss function to a convex function.

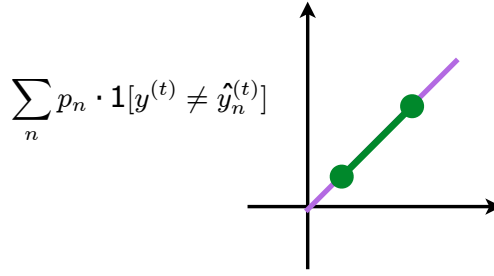


Figure 9: Expectation over zero-one loss function is now convex

### b. Convexification by surrogate loss

Instead of dealing with non-convex loss functions, we can design a convex function that upper-bounds the original loss function and solve the optimization problem using convex optimization infrastructure. This idea is called convexification by surrogate loss. We use the weighed majority algorithm again as an example. The original non-convex zero-one loss function is:

$$l^{(t)} = \mathbf{1}[\hat{y}^{(t)} \neq y^{(t)}]$$

We can design a surrogate loss function that upper bounds this loss function to approximate its behavior while making the problem convex. To do so, we need to find a theoretical upper bound for this zero-one loss:

$$\begin{aligned} l^{(t)} &= \mathbf{1}[\hat{y}^{(t)} \neq y^{(t)}] \\ &= \mathbf{1}[-y^{(t)} \hat{y}^{(t)} > 0] && y \text{ is either -1 or 1} \\ &= \mathbf{1}[-y^{(t)} \langle \mathbf{w}^{(t-1)}, \mathbf{x}^{(t)} \rangle > 0] && \text{substitute (2.3.4)} \\ &= \mathbf{1}[1 - y^{(t)} \langle \mathbf{w}^{(t-1)}, \mathbf{x}^{(t)} \rangle > 1] && \text{shift to the right by 1} \end{aligned}$$

We can upper bound the original loss function by using a hinge function, as visualized in Fig. 10:

$$\begin{aligned} l^{(t)} &= \mathbf{1}[1 - y^{(t)} \langle \mathbf{w}^{(t-1)}, \mathbf{x}^{(t)} \rangle > 1] \\ &\leq \max \left[ 0, 1 - y^{(t)} \langle \mathbf{w}^{(t-1)}, \mathbf{x}^{(t)} \rangle \right] && \text{a convex hinge function } \tilde{l}^{(t)} \end{aligned}$$

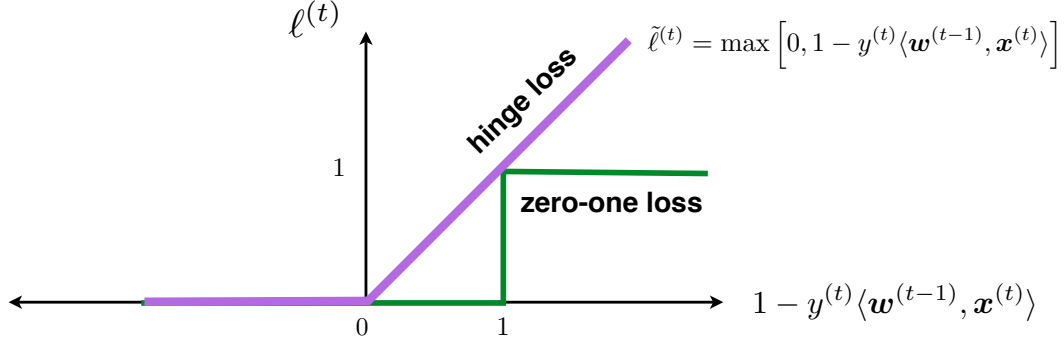


Figure 10: The hinge function acts as a surrogate loss for the zero-one loss function by upper-bounding it

By optimizing the hinge loss  $\tilde{\ell}^{(t)}$ , we have a slightly different behavior than the original zero-one loss especially in the  $[0, 1]$  domain where the hinge loss produces non-zero penalty while the original loss function does not. This is a trade-off where we make the optimization easier to solve (using an OCO solver) at the cost of an approximation.

### 3 Follow the Leader

After reviewing the theory part of OCO, we introduce an algorithm called *Follow the Leader* (FTL) known as ‘Fictitious Play’ (Brown 1951 [1]). The idea is that the learner select the best choice of the game so far. The FTL pseudocode is listed below.

---

**Algorithm 4** Follow the leader algorithm

---

```

1: function FOLLOW THE LEADER
2:   for  $t = 1, 2, \dots, T$  do
3:      $\mathbf{w}^{(t)} = \arg \min_{\mathbf{w} \in W} \sum_{i=1}^{t-1} f^{(i)}(\mathbf{w})$  ▷ Parameter chosen
4:     RECEIVE  $(f^{(t)} : W \rightarrow \mathbb{R})$  ▷ Receive loss function
5:   end for
6: end function

```

---

The prediction rule is that the learner is going to choose the parameter  $\mathbf{w}$  such that it minimizes the cumulative loss up to now. This rule can be inserted into the online convex optimizer as mentioned in Fig. 3.

Now we derive a general regret bound for the FTL algorithm. First, we use one-step look ahead cheater to get the regret upper bound. Then we apply algebraic tricks to get the final form.

Recall the definition of cumulative regret:

$$R(\mathbf{u}) = \sum_t [f^{(t)}(\mathbf{w}^{(t)}) - f^{(t)}(\mathbf{u})] \quad (2)$$

where  $\mathbf{u}$  is any hypothesis or parameter. Therefore, the cumulative regret is the sum of all the

differences between the loss of your online convex optimizer and the loss of hypothesis  $\mathbf{u}$ .

Because hypothesis  $\mathbf{u}$  has infinite possibilities, it is very hard to reason about  $\mathbf{u}$ . Therefore, we introduce “the one-step look ahead cheater” to upper bound the regret so that we can analyze.

$$R(\mathbf{u}) = \sum_t [f^{(t)}(\mathbf{w}^{(t)}) - f^{(t)}(\mathbf{u})] \leq \sum_t [f^{(t)}(\mathbf{w}^{(t)}) - f^{(t)}(\mathbf{w}^{(t+1)})] \quad (3)$$

where  $\mathbf{w}^{(t+1)}$  is called “one-step look ahead cheater”. It means that the prediction rule at step  $t$  that has access to parameters of the next iteration  $t + 1$ . We don’t have access to the parameters of the future. However, we can use theoretical values to enable the analysis. The “one-step look ahead cheater” gives us the access to future observations/loss that enables the computation of the best  $\mathbf{w}$ .

To prove the inequality (3), by subtracting  $f^{(t)}(\mathbf{w}^{(t)})$  from both sides, we only need to prove that:

$$\sum_{t=1}^T f^{(t)}(\mathbf{w}^{(t+1)}) \leq \sum_{t=1}^T f^{(t)}(\mathbf{u}) \quad (4)$$

If it holds true, it means that the loss of a series of one step cheaters is not larger than the loss of any other single parameter  $\mathbf{u}$ . Therefore, our claim is:

$$\sum_{t=1}^T f^{(t)}(\mathbf{w}^{(t+1)}) \leq \sum_{t=1}^T f^{(t)}(\mathbf{u}) \quad \forall \mathbf{u} \quad (5)$$

Proof. By induction method, we first assume our claim is true for  $T - 1$  and will prove it holds true for  $T$ :

$$\sum_{t=1}^{T-1} f^{(t)}(\mathbf{w}^{(t+1)}) \leq \sum_{t=1}^{T-1} f^{(t)}(\mathbf{u}) \quad (6)$$

Then we add  $f^{(T)}(\mathbf{w}^{(T+1)})$  to both sides of inequality (6), where  $f^{(T)}(\mathbf{w}^{(T+1)})$  is the loss of the one step look ahead cheater for the next time step, we can get:

$$\sum_{t=1}^T f^{(t)}(\mathbf{w}^{(t+1)}) \leq \sum_{t=1}^{T-1} f^{(t)}(\mathbf{u}) + f^{(T)}(\mathbf{w}^{(T+1)}) \quad (7)$$

Since the inequality (7) holds true for all  $\mathbf{u}$ , if we let  $\mathbf{u} = \mathbf{w}^{(T+1)}$ , we can get:

$$\sum_{t=1}^T f^{(t)}(\mathbf{w}^{(t+1)}) \leq \sum_{t=1}^{T-1} f^{(t)}(\mathbf{w}^{(T+1)}) + f^{(T)}(\mathbf{w}^{(T+1)}) \quad (8)$$

$$\sum_{t=1}^T f^{(t)}(\mathbf{w}^{(t+1)}) \leq \sum_{t=1}^T f^{(t)}(\mathbf{w}^{(T+1)}) \quad (9)$$



It means that the loss of always using one step look ahead cheater is the lower bound of the loss of using a single cheater at the end of the sequence for all time steps. Recall the minimizer by definition,

$$\mathbf{w}^{(T+1)} = \arg \min_{\mathbf{u} \in S} \sum_{t=1}^T f^{(t)}(\mathbf{u}) \quad (10)$$

This means the right hand side of the inequality (9) can be upper bounded by  $\sum_{t=1}^T f^{(t)}(\mathbf{u})$ . Therefore, we have

$$\sum_{t=1}^T f^{(t)}(\mathbf{w}^{(t+1)}) \leq \sum_{t=1}^T f^{(t)}(\mathbf{u}) \quad (11)$$

Therefore, it proves our claim and the inequality (3).

## References

- [1] U. Berger. Brown’s original fictitious play. *Journal of Economic Theory*, 135(1):572–578, 2007.
- [2] Wikipedia contributors. Lipschitz continuity — Wikipedia, the free encyclopedia, 2021. [Online; accessed 17-Feb-2021].